

Configuration Manual

MSc Research Project M.Sc. Cybersecurity

Jerry Devassy Student ID: x19220961

School of Computing National College of Ireland

Supervisor: Liam McCabe

National College of Ireland



Year: 2021-2022

9

MSc Project Submission Sheet

School of Computing

Student Name: Jerry Jockey Devassy

Student ID: X19220961

Programme: M.Sc. Cybersecurity

Module:

Lecturer: Liam McCabe Submission Due

Date: 16/12/2021

- **Project Title:** Detection of Application Layer DDoS Attack by using logistic Regression
- Word Count: 1241 Page Count:

MSc Research Project

I hereby certify that the information contained in this (my submission) is information pertaining to research I conducted for this project. All information other than my own contribution will be fully referenced and listed in the relevant bibliography section at the rear of the project.

<u>ALL</u> internet material must be referenced in the bibliography section. Students are required to use the Referencing Standard specified in the report template. To use other author's written or electronic work is illegal (plagiarism) and may result in disciplinary action.

Signature: Jerry Jockey Devassy

Date: 16th December 2021

PLEASE READ THE FOLLOWING INSTRUCTIONS AND CHECKLIST

Attach a completed copy of this sheet to each project (including multiple copies)	
Attach a Moodle submission receipt of the online project	
submission , to each project (including multiple copies).	
You must ensure that you retain a HARD COPY of the project,	
both for your own reference and in case a project is lost or mislaid. It is	
not sufficient to keep a copy on computer.	

Assignments that are submitted to the Programme Coordinator Office must be placed into the assignment box located outside the office.

Office Use Only	
Signature:	
Date:	
Penalty Applied (if applicable):	

Configuration Manual

Jerry Jockey Devassy Student ID: X19220961

1 Introduction

This study requires the usage of a dataset that includes both DDoS attack traffic and benign traffic characteristics. All irrelevant data present in the dataset such as NaaN, Negative and null values are deleted as part of the data pre-processing procedure. It is also important to do feature selection, encoding the class variables, forming a subset of datasets, developing the model, and evaluating the model. The setup manual's objective is to help users in installing this research project code on their system so that they may assess the study or make changes to match their individual requirements. The prerequisites and environment set up section povides comprehensive information for constructing a project setting as well as a list of need for reproducing the study results. The Code Execution section includes the completed created code as well as the parameters for modifying various aspects of the project.

2 Requirement

To implement the project, a system that supports all the essential tools and related settings is required, so that the implemented project runs successfully with all of its supported dependencies. If the system supports all of the essential tools, there will be no issues running the project resulting in a better outcome.

2.1 System Requirement

The machine learning process consumes very high resources on the host system. As a result, the hardware configuration on the system in which the project is downloaded must be capable of doing such duties. The system's minimal requirements are as follows:

For Windows:

- CPU: Intel i3 5th Gen and above
- RAM: 8GB DDR4 and above
- Storage: 20 GB free space in HDD

For iOS:

- Processor: 1.1GHz dual core Intel Core i3
- Memory: 8GB LPDDR4X onboard memory
- Storage: 128GB PCIe-based SSD

2.2 Machine Requirement

- MS Excel for analysing the dataset
- Proper working Internet Connection
- Web Browsers Safari / Google Chrome

2.3 Software Requirements

• Anaconda Navigator- Jupyter Notebook("Anaconda | Choose Your Anaconda IDE Adventure," 2021)



Figure 1: Anaconda Navigator- Jupyter Notebook

• Python 3- It is an open-sourced software available to download. We used python 3 as it has better features available as compared to python 2("Download Python," 2021).



Figure 2: Python Download

3 Dataset

The CICIDS 2017 dataset was used in this research. The dataset is a CSV file having normal traffic and four different DDoS attack traffic captured from the web traffic. This dataset is open sourced and is available to download from the CIC Website("IDS 2017 | Datasets | Research | Canadian Institute for Cybersecurity | UNB," 2017).

Figure 3 shows the license of the dataset and how to download the dataset. however, it will require a strong internet connection as the size of the dataset is 2GB. Figure 2 shows the sample of dataset in the excel sheet.

License

The CICIDS2017 dataset consists of labeled network flows, including full packet payloads in pcap format, the corresponding profiles and the labeled flows (GeneratedLabelledFlows.zip) and CSV files for machine and deep learning purpose (MachineLearningCSV.zip) are publicly available for researchers. If you are using our dataset, you should cite our related paper which outlining the details of the dataset and its underlying principles:

 Iman Sharafaldin, Arash Habibi Lashkari, and Ali A. Ghorbani, "Toward Generating a New Intrusion Detection Dataset and Intrusion Traffic Characterization", 4th International Conference on Information Systems Security and Privacy (ICISSP), Portugal, January 2018

Download this dataset



Figure 4: Preview of Dataset

4 Required Packages and Imports

The model was implemented on jupyter notebook and python code was used. The Python code includes packages and imports that are listed below:

- Matplotlib 3.1.2
- Numpy 1.18.0
- scikit-learn 0.22
- sklearn
- pandas
- mpl_toolkits
- Seaborn
- train_test_split
- Sklearn.metrics
- Logisticregression

5 Execution of Code

Below are the few steps to run the code

- 1. Download the zip file with datasets and python files
- 2. Unzip the files into a folder
- 3. Open Jupyter notebook through anaconda navigator GUI
- 4. Open .ipynb file in Jupyter (DDoS_Detection_thesis)
- 5. Run the code, either step by step or whole code at the same time.

6 Evaluation of Code

There are multiple sections in which the whole model is implemented. Firstly, the whole data is sent to feature selection in which correlation coefficient method is used to select the best possible features and later these features is sent to logistic regression model to train and test the model and to classify if the traffic is an attack or a normal traffic.

6.1 Feature Selection and Data Processing

For feature selection, Correlation Coefficient method is used to select the best features that are dependent to each other. In Figure 5, function for plotting of correlation matrix is written but as the size of the dataset is large the plotting shows the graph but has no clear values in it. Figure 5 shows the function of correlation matrix in which it is implemented.

```
In [211]: # Correlation matrix
def plotCorrelationMatrix(df, graphWidth):
    #filename = df.dataframeName
    df = df.dropna('columns') # drop columns with NaN
    df = df[[col for col in df if df[col].nunique() > 1]] # keep columns where there are more than 1 unique va
    if df.shape[1] < 2:
        print(f'No correlation plots shown: The number of non-NaN or constant columns ({df.shape[1]}) is less
        return
        corr = df.corr()
    plt.figure(num=None, figsize=(graphWidth, graphWidth), dpi=80, facecolor='w', edgecolor='k')
        corrMat = plt.matshow(corr, fignum = 1)
        plt.xticks(range(len(corr.columns)), corr.columns, rotation=90)
        plt.yticks(range(len(corr.columns)), corr.columns)
        plt.gca().xaxis.tick_bottom()
        plt.title(f'Correlation Matrix for ', fontsize=15)
        plt.show()</pre>
```

Figure 5: Function for Correlation matrix

The preceding code is used to generate the correlation function, which is called by giving the data frame and the graph width in it. Figure 6 shows how the function is called and the graph width that is mentioned in it.

Figure 6: Plotting the Correlation function

The Output of the plotting gives the features that are strongly in relation with the target variable (variable that states whether it is legit traffic or attack traffic). Since the data was too large, the plotting matrix did not show proper data in it and therefore in Figure 7, A function has been implemented to retrieve the features which are highly correlated with each other against the target variable (Label) are called through the function itself which lists down the features in the output.

Figure 7: List of Features

6.2 Using Logistic Regression to classify the Data

Since the dataset was huge (count of traffic in dataset was around 6lakhs), we had to create the subset of dataset from the original one which was performed in a random manner and imported the same in the code which is listed in Figure 8.



Figure 8: Data fed into Logistic Regression

In Figure 9, StandardScaler functions is used because in StandardScaler the mean is removed, and each feature/variable is scaled to unit variance. The data is divided into X – independent variables and y – target variables, which are subsequently utilized by the logistic regression classifier to train and test the model.

```
In [71]:
          # Feature Scaling - which might help to improve the performance of the algorithm
          #from sklearn.preprocessing import StandardScaler
          scalar = StandardScaler()
          X_train = scalar.fit_transform(X_train)
          X_test = scalar.transform(X_test)
In [461]: # Fitting the Logistic Regression into the Training set
          #from sklearn.linear_model import LogisticRegression
          #model_lg = LogisticRegression()
          model_lg = LogisticRegression(solver='lbfgs', max_iter=5000)
          model_lg.fit(X_train,Y_train)
          # Predicting the test set results
          Y_test_pred = model_lg.predict(X_test)
          Y_test_pred
Out[461]: array([1, 0, 0, ..., 0, 1, 0], dtype=int64)
In [463]: model_lg.score(X_test, Y_test)
Out[463]: 0.82833333333333334
```



The model's output is assessed in terms of an accuracy score, which is computed using the confusion matrix. Figure 10 clearly depicts the confusion matrix and accuracy obtained while training the model using logistic regression.

```
In [464]: # Making the Confusion Matrix
#from sklearn.metrics import accuracy_score, confusion_matrix
confusion_matrix(Y_test, Y_test_pred)
Out[464]: array([[1341, 185],
      [ 330, 1144]], dtype=int64)
In [465]: accuracy_score (Y_test,Y_test_pred)
Out[465]: 0.828333333334
```

Figure 10: Accuracy and Confusion Matrix

7 References

- Anaconda | Choose Your Anaconda IDE Adventure: Jupyter, JupyterLab,... [WWW Document], 2021 . Anaconda. URL https://www.anaconda.com/blog/choose-your-anaconda-ide-adventure-jupyter-jupyterlab-or-apache-zeppelin (accessed 12.14.21).
- Download Python [WWW Document], 2021 . Python.org. URL https://www.python.org/downloads/ (accessed 12.14.21).
- IDS 2017 | Datasets | Research | Canadian Institute for Cybersecurity | UNB [WWW Document], 2021. URL https://www.unb.ca/cic/datasets/ids-2017.html (accessed 12.13.21).